

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-319768

(43) 公開日 平成9年(1997)12月12日

(51) Int.Cl. ⁸	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 17/30			G 0 6 F 15/401	3 2 0 A
15/18	5 6 0		15/18	5 6 0 J

審査請求 未請求 請求項の数2 F D (全 5 頁)

(21) 出願番号 特願平8-157723

(22) 出願日 平成8年(1996)5月29日

(71) 出願人 000000295

沖電気工業株式会社

東京都港区虎ノ門1丁目7番12号

(72) 発明者 福本 淳一

東京都港区虎ノ門1丁目7番12号 沖電気
工業株式会社内

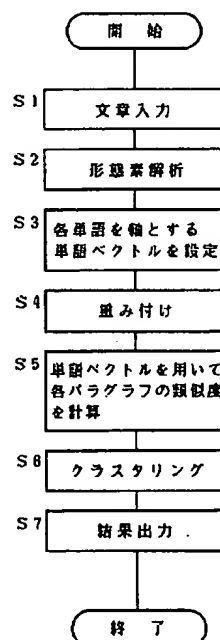
(74) 代理人 弁理士 佐藤 幸男 (外1名)

(54) 【発明の名称】 要点抽出方法

(57) 【要約】

【課題】 文章中の各文中の語句の参照関係や文章の構造情報等を用いることなく、文章中の重要部分を抽出する。

【解決手段】 要点抽出対象の文章が入力されると(ステップS1)、形態素解析を行って(ステップS2)、単語を抽出する。次いで、各単語を軸とする単語ベクトルを設定すると共に(ステップS3)、各単語に重み付けの値を付与する(ステップS4)。更に、各パラグラフの類似度を単語ベクトルの値を用いて計算する(ステップS5)。このパラグラフの類似度からクラスタを生成し(ステップS6)、その結果を出力する(ステップS7)。



本発明方法のフローチャート

1

【特許請求の範囲】

【請求項1】 自然言語で記述された文章中の各パラグラフを、当該文章中出现する各単語を軸とし、各軸に対して、各々のパラグラフ中出现する回数を対応させた単語ベクトルで表し、

各パラグラフの単語ベクトルの類似度を計算し、前記単語ベクトルの最も類似度の高いパラグラフを、前記文章における最も重要な部分として出力することを特徴とする要点抽出方法。

【請求項2】 請求項1記載の要点抽出方法において、任意の単語は、予め決められた基準により重み付けされることを特徴とする要点抽出方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、自然言語で記述された文章から重要部分を抽出する要点抽出方法に関するものである。

【0002】

【従来の技術】従来の自然言語で記述された文章から重要部分を抽出する方法においては、文章中の各文の形態素解析、構文解析を行い、各文中の語句の参照関係に基づき、他の文から最も多く参照されている文を重要文としたり、接続語句等を用いた文章の構造情報を用いて重要部分を判定するといったヒューリスティクスに基づく手法が用いられていた。

【0003】

【発明が解決しようとする課題】しかしながら、一般に、文章中の各文中の語句の参照関係の解析のためには、その文章中で用いられる語句に関する上位・下位の知識等の多くの知識が必要であり、参照関係の解析そのものが困難であるといった問題があった。また、接続語句等を用いて得られた文章の構造情報から重要部分を判定する方法もあるが、そのためには文章の構造を解析するための接続語句等の情報を多く登録しておかなければならないといった問題があった。

【0004】このような点から、文章中の各文中の語句の参照関係や文章の構造情報等を用いることなく、文章中の重要部分を抽出することのできる要点抽出方法の実現が望まれていた。

【0005】

【課題を解決するための手段】本発明は、前述の課題を解決するため次の構成を採用する。

〈請求項1の構成〉自然言語で記述された文章中の各パラグラフを、その文章中出现する各単語を軸とし、各軸に対して、各々のパラグラフ中出现する回数を対応させた単語ベクトルで表し、各パラグラフの単語ベクトルの類似度を計算し、単語ベクトルの最も類似度の高いパラグラフを、文章における最も重要な部分として出力することを特徴とする要点抽出方法である。

【0006】〈請求項1の説明〉請求項1の発明は、文

2

章中、重要な部分は、繰り返し述べられていることが多いという点に着目し、同様な事柄が述べられているパラグラフをその文章の要点であると判定するようにしたものである。ここで、パラグラフとしては、文章中の章単位、文単位等、任意の区切りであってもよい。また、軸とする単語は、文章中出现する全ての単語を対象とするが、適宜選択するようにしてもよい。

【0007】このように、文章中の単語のみの情報を用いて、それらの情報の統計的処理により、文章中の重要なパラグラフを抽出するようにしているため、文章中の各文中の語句の参照関係や文章の構造情報等を用いることなく、容易かつ正確に重要部分を抽出することが可能となる。

【0008】〈請求項2の構成〉請求項1記載の要点抽出方法において、任意の単語は、予め決められた基準により重み付けされることを特徴とする要点抽出方法である。

【0009】〈請求項2の説明〉予め決められた基準とは、例えば単語の品詞情報に基づく基準である。即ち、助詞、助動詞等の付属語は文章中に多く出現するが、これらの付属語情報は文章の重要度の判定には必要ないため、低い重み付けの値を設定する。一方、文章中の各単語のうち、名詞や動詞等の自立語は文章中で重要部分を判定するために必要であるため、高い重み付けの値を設定する。これにより、各パラグラフに対して単語の重要度を考慮した単語ベクトルが設定される。また、この重み付けの基準は、単語の品詞情報だけでなく、これ以外にも、ユーザによって特定の単語を指定するといったように、適宜選択が可能である。

【0010】

【発明の実施の形態】以下、本発明の実施の形態を図面を用いて詳細に説明する。図1は本発明の要点抽出方法を示すフローチャートであるが、この説明に先立ち、本発明の要点抽出方法を実現するための要点抽出装置を説明する。

【0011】図2は、その要点抽出装置を示す構成図である。図の装置は、入力部1、形態素解析処理部2、重要部分抽出部3、出力部4、重み付け処理部5からなる。

【0012】要点抽出装置は、マイクロコンピュータで構成され、入力部1は、例えば入力インタフェースやキーボードといった解析対象文の入力を行う部分である。また、形態素解析処理部2は、入力部1に入力された文を各単語に分割する処理を行う機能を有している。

【0013】重要部分抽出部3は、形態素解析された単語情報と各単語に付与された重み付け情報を用いて重要部分であるパラグラフを抽出する機能を有している。即ち、この重要部分抽出部3は、重要な部分は、文章中で繰り返し出現することが多いという点に着目し、同様な事柄が述べられているパラグラフ、つまり、共通してい

3

る単語が最も多いパラグラフを、その文章における要点が記述されているパラグラフとして出力するようにしたものである。

【0014】重み付け処理部5は、形態素解析された単語情報に対して、重み付けの計算を行う機能を有している。また、出力部4は、例えば、表示装置や印刷装置といった出力部であり、重要部分抽出部3で抽出された重要部分の出力を行う機能を有している。尚、上記の形態素解析処理部2～重み付け処理部5は、各機能を実現するプログラムと、これを実行するマイクロコンピュータにおける中央処理装置やメモリといった制御部により構成されている。

【0015】次にこのように構成された要点抽出装置を用いた要点抽出方法を図1に沿って説明する。まず、ユーザは、要点抽出を行う文を入力部1に入力する（ステップS1）。これにより、形態素解析処理部2は文章中の各パラグラフにおける単語の認識処理を行う（ステップS2）。尚、この形態素解析処理については既知の処理であるため、ここでの説明は省略する。

【0016】次に、重要部分抽出部3は、文章中の各パラグラフ中に存在する全ての単語情報に対して、各単語を軸とする単語ベクトルを設定する（ステップS3）。図3は、各パラグラフの単語ベクトルの説明図である。この例は、三つの単語で四つのパラグラフの場合を示しており、図中、軸6、7、8がそれぞれ、単語1、単語2、単語3を示し、9～12が、各パラグラフに対する単語ベクトルを示している。

【0017】また、重要部分の抽出のためには、単語の重要度が異なるため、重み付け処理部5は、それを表す重み付けの値を、抽出された各単語に対して付与する（図1におけるステップS4）。この重み付けの値を与える方法としては、例えば、文章中の各単語の品詞情報を用いる方法がある。これは、文章中の各単語のうち、助詞、助動詞等の付属語は文章中に多く出現するが、これらの付属語情報は文章の重要度の判定には必要ないため、低い重み付けの値を設定する。一方、文章中の各単語のうち、名詞や動詞等の自立語は文章中で重要部分を判定するために必要であるため、高い重み付けの値を設定する。そして、文章中の各単語の品詞情報から設定された重み付けの値を文章中の各パラグラフの単語ベクトルに対して掛け合わせる。これにより、各パラグラフに対して単語の重要度を考慮した単語ベクトルが設定される。尚、このような重み付けの値の付与は、これ以外にも、ユーザが、重み付け処理部5に対して特定の単語を任意の値を指定できるよう構成してもよい。

【0018】次に、文章中の各パラグラフに対して、類似度の計算を、上記の重み付けされた単語ベクトルを用いて行う（ステップS5）。即ち、単語ベクトルの値が類似しているパラグラフを類似度の高いパラグラフであると判断する。尚、このような類似度の計算方法とし

(3)

特開平9-319768

4

て、例えば「G. Salton: Automatic Text Processing, Addison-Wesley Publishing Company (1989)」Chapter 10等示されている方法を用いることができる。

【0019】全てのパラグラフの類似度の計算が済むと、その類似度を用いて各パラグラフについてクラスタリングを行う（ステップS6）。そして、クラスタリングが終了すると、その結果を、出力部4に出力し（ステップS7）、要点抽出処理を終了する。

【0020】次に、上記の動作を更に詳細に説明する。文章中の各パラグラフに対して、重み付けされた単語ベクトルに基づきその類似度の計算による重要部分の抽出の手法を以下に示す。

【0021】図4は、要点抽出処理における演算式の説明図である。パラグラフの類似度計算の方法としては、パラグラフの単語ベクトルを $X = (x_1, x_2, \dots, x_t)$ 、 $Y = (y_1, y_2, \dots, y_t)$ とした場合、単語ベクトル X 、 Y の類似度は、図中の式(1)で表される。

【0022】また、単語1、単語2の単語ベクトルをそれぞれ $W1 = (w_{11}, w_{12}, \dots, w_{1t})$ 、 $W2 = (w_{21}, w_{22}, \dots, w_{2t})$ とし、単語の重み付けベクトルを $A = (a_1, a_2, \dots, a_t)$ としたとき、単語ベクトル $W1$ 、 $W2$ の類似度 S_{12} は、図中の式(2)で計算される。尚、このとき、 t は単語の種類数である。

【0023】以上の類似度の計算を文章中の全てのパラグラフ(1, ..., n)について計算した結果は、図中の配列(3)のように示される。尚、ここで、 S_{ij} はパラグラフ i とパラグラフ j の類似度を計算した値であるとする。但し、 $S_{ii} = 0$ であるとする。例えば、配列(4)は、パラグラフ1, 2, 3, 4からなる文章について得られた配列を示している。

【0024】次に、以上のようにして得られた各パラグラフ間の類似度の値を用いてパラグラフのクラスタリングを行う。このクラスタリングの方法としては、上述した「G. Salton: Automatic Text Processing, Addison-Wesley Publishing Company (1989)」Chapter 10において示されている方法を用いることが可能である。

【0025】図5は、クラスタリングの一例を示す図である。この例は、上記図4の配列(4)で示したパラグラフのクラスタリングを示している。即ち、配列(4)において、類似度の最も高いのは、 S_{14} および S_{41} の0.9である。従って、文章中、パラグラフ1とパラグラフ4とが最も類似度の高いパラグラフであるため、これらのパラグラフからクラスタを生成する。

【0026】次に、パラグラフ1またはパラグラフ4のどちらかのパラグラフと類似度の高いパラグラフを抽出する。ここで、パラグラフ1と最も類似度の高いパラグラフの値は、パラグラフ2との0.7であり、また、パ

10

20

30

40

50

5

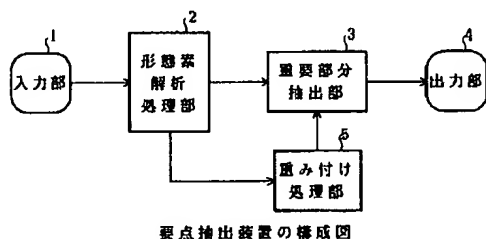
ラグラフ4と最も類似度の高い値は、パラグラフ2との0.5である。従って、パラグラフ1とパラグラフ4とで生成したクラスとパラグラフ2とで上位のクラスを生成する。

【0027】以下、同様に、生成したクラスのうち、いずれかのパラグラフと最も類似度の高いパラグラフで、更に上位のクラスを生成する。ここでは、パラグラフが4個であるため、残りのパラグラフ3によって上位のクラスが生成される。尚、クラス生成を、いずれかのパラグラフの一方との比較ではなく、二つのパラグラフの合成ベクトルとの比較によって行うようにしてもよい。

【0028】そして、このようなクラスタリングの結果が出力部4から出力される。これにより、ユーザは、パラグラフ1とパラグラフ4とが最も重要なパラグラフであることを知ることができる。即ち、文章中、パラグラフ1とパラグラフ4とで同様の事柄が最も多く述べられているため、これらのパラグラフで文章の要点が記されていることが分かる。

【0029】以上のように、上記具体例では、文章中の各パラグラフの単語のみの情報を用い、それらの情報の統計的処理により、パラグラフ間の類似度を計算して、

【図2】



【図4】

(1)	$\sum_{i=1}^t x_i \cdot y_i$
(2)	$S_{12} = \sum_{i=1}^t a_i w_{1i} \cdot a_i w_{2i}$
(3)	$\begin{pmatrix} S_{11} & S_{12} & \dots & S_{1n} \\ S_{21} & S_{22} & \dots & S_{2n} \\ \dots & \dots & \dots & \dots \\ S_{n1} & S_{n2} & \dots & S_{nn} \end{pmatrix}$
(4)	$\begin{pmatrix} S_{11} & S_{12} & S_{13} & S_{14} \\ S_{21} & S_{22} & S_{23} & S_{24} \\ S_{31} & S_{32} & S_{33} & S_{34} \\ S_{41} & S_{42} & S_{43} & S_{44} \end{pmatrix} = \begin{pmatrix} 0 & 0.7 & 0.6 & 0.9 \\ 0.7 & 0 & 0.5 & 0.5 \\ 0.6 & 0.5 & 0 & 0.4 \\ 0.9 & 0.5 & 0.4 & 0 \end{pmatrix}$

演算式の説明図

6

文章中の重要部分の判定を行うようにしたので、文章中の各文中の語句の参照関係や文章の構造情報等を用いることなく、容易にかつ正確に要点を抽出することができる。

【図面の簡単な説明】

【図1】本発明の要点抽出方法を示すフローチャートである。

【図2】本発明の要点抽出方法を実現するための要点抽出装置の構成図である。

10 【図3】本発明の要点抽出方法における各パラグラフと単語ベクトルの説明図である。

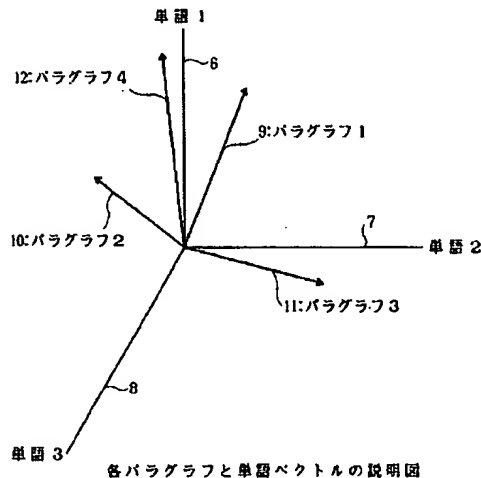
【図4】本発明の要点抽出方法における演算式の説明図である。

【図5】本発明の要点抽出方法におけるクラスタリングの一例を示す図である。

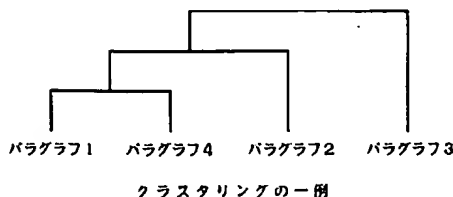
【符号の説明】

- 1 入力部
- 2 形態素解析処理部
- 3 重要部分抽出部
- 4 出力部
- 5 重み付け処理部

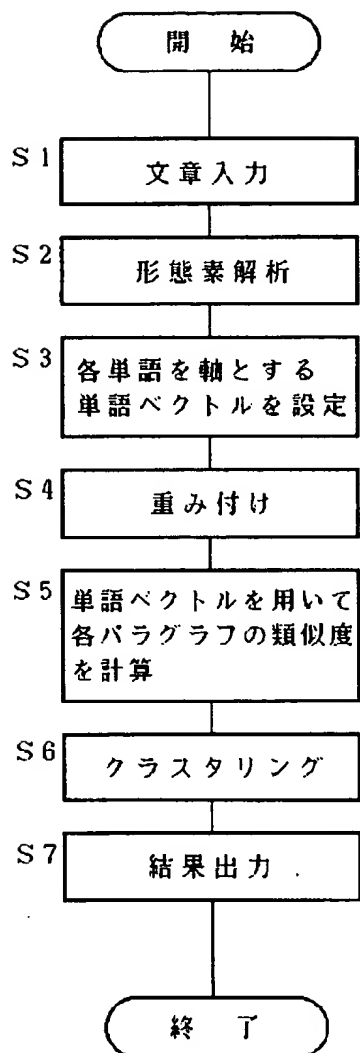
【図3】



【図5】



【図1】



本発明方法のフローチャート